

# Northeastern University

## Solving Interpretable Kernel Dimension Reduction

Chieh Wu, Jared Miller, Yale Chang, Mario Sznaier, Jennifer G. Dy  
Dept .of Electrical and Computer Engineering, Northeastern University

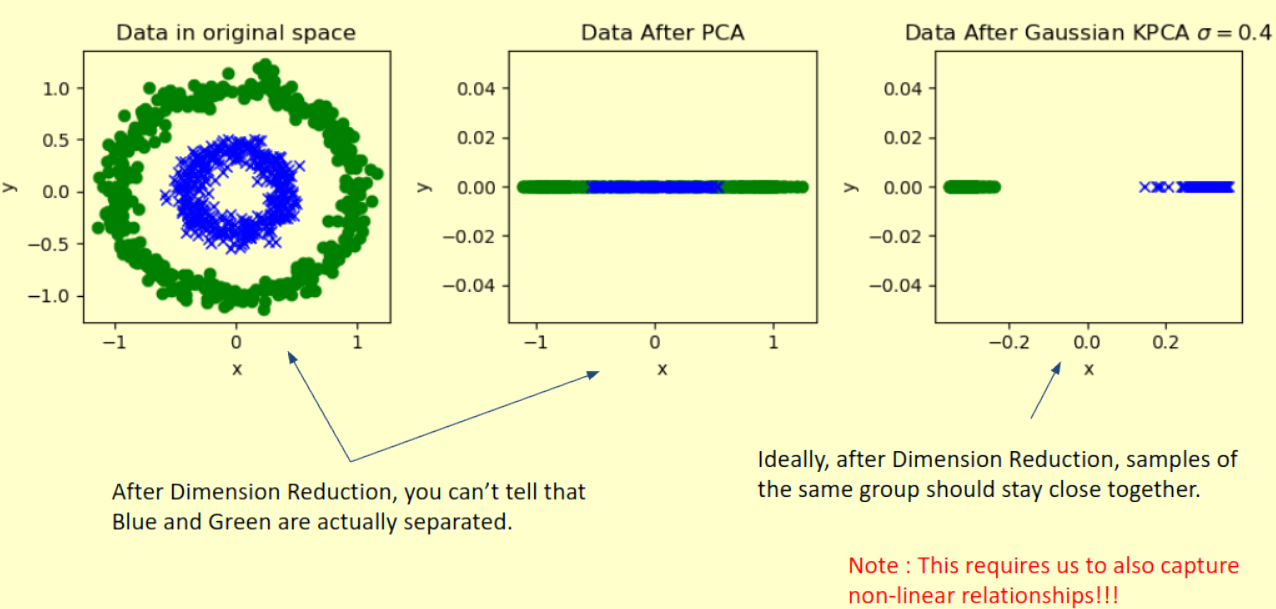
### What is IKDR?

Principal Component Analysis (PCA) is the most commonly used Dimension Reduction (DR) technique. It is also an **interpretable** way to reduce the dimension.

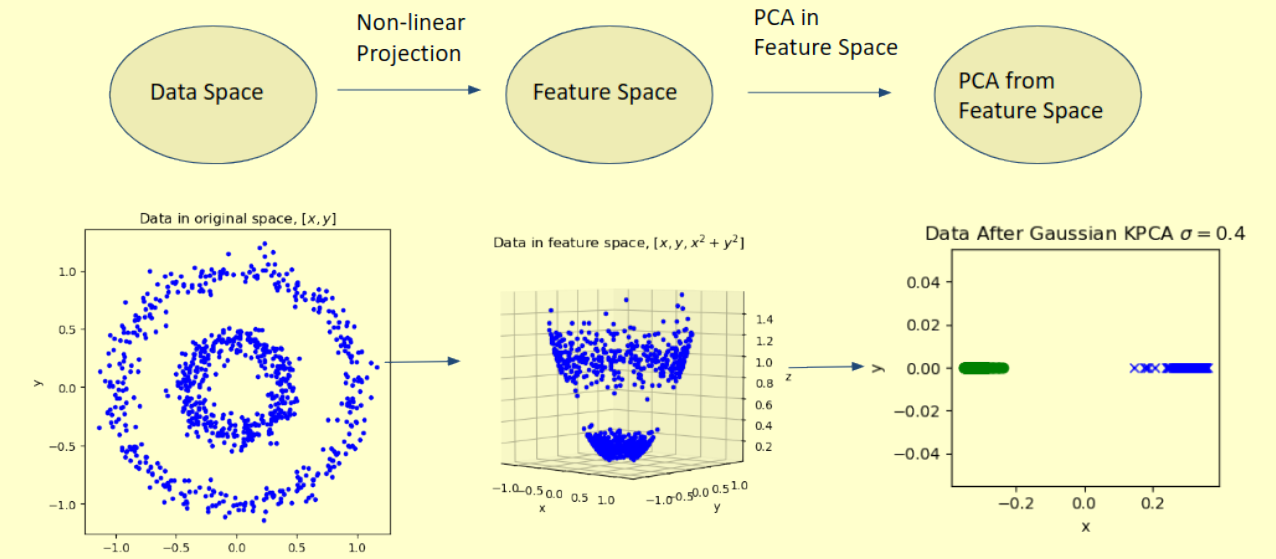
We know exactly how the new features relate to the original features.

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} x_1 w_{11} + x_2 w_{12} + x_3 w_{13} \\ x_1 w_{21} + x_2 w_{22} + x_3 w_{23} \end{bmatrix}$$

But PCA cannot capture **nonlinear Relationships**.



**KPCA captures nonlinear Relationships but not interpretable.**



**KPCA is very powerful, but .....**

Problem 1: It does not use labels to guide the dimension reduction.

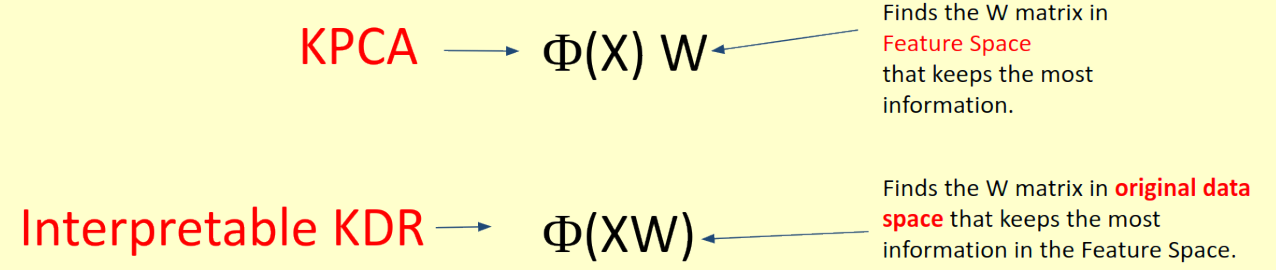
Problem 2: Since KPCA is PCA in the feature space, it's not obvious what they mean.

Here is the Gaussian Kernel feature map:

$$\phi(x) = e^{-x^2/2\sigma^2} \left[ 1, \sqrt{\frac{1}{1!\sigma^2}}x, \sqrt{\frac{1}{2! \sigma^4}}x^2, \sqrt{\frac{1}{3! \sigma^6}}x^3, \dots \right]^T$$

**Interpretable Kernel Dimension Reduction (IKDR) solves both problems...**

How IKDR produce interpretable results.



$$\max_W HSIC(XW, Y) \quad \text{s.t.} \quad W^T W = I$$

HSIC(X,Y) measures the non-linear dependence between X and Y in Feature Space.

Although this make the solution interpretable, it is very difficult to solve.

In general, many IKDR problems have a common objective.

$$\max_W \sum_{i,j} \Gamma_{i,j} K_{XW_i,j} \quad \text{s.t.} \quad W^T W = I \quad (1)$$

Symmetric Positive Definite Matrix      Kernel Matrix on XW      Grassmann Manifold Constraint.

### Where is IKDR used?

**Supervised Dimension Reduction for Classification**

$$\max_W HSIC(XW, Y) \quad \text{s.t.} \quad W^T W = I$$

**Unsupervised Dimension Reduction for Clustering**

$$\max_{W,Y} HSIC(XW, Y) \quad \text{s.t.} \quad W^T W = I$$

**Semi-supervised Dimension Reduction for Clustering Using Multiple Expert Sources**

$$\max_{W,Y} \text{Tr}(Y^T \mathcal{L}_W Y) + \mu \text{Tr}(K_{XW} H K_{\hat{Y}} H)$$

$$\text{s.t.} \quad \mathcal{L}_W = D^{-\frac{1}{2}} K_{XW} D^{-\frac{1}{2}} W^T W = I, Y^T Y = I$$

**Alternative Clustering via Dimension Reduction**

$$\max_{W,Y} \text{Tr}(K_{XW} H K_Y H) - \mu \text{Tr}(K_{XW} H K_{\hat{Y}} H)$$

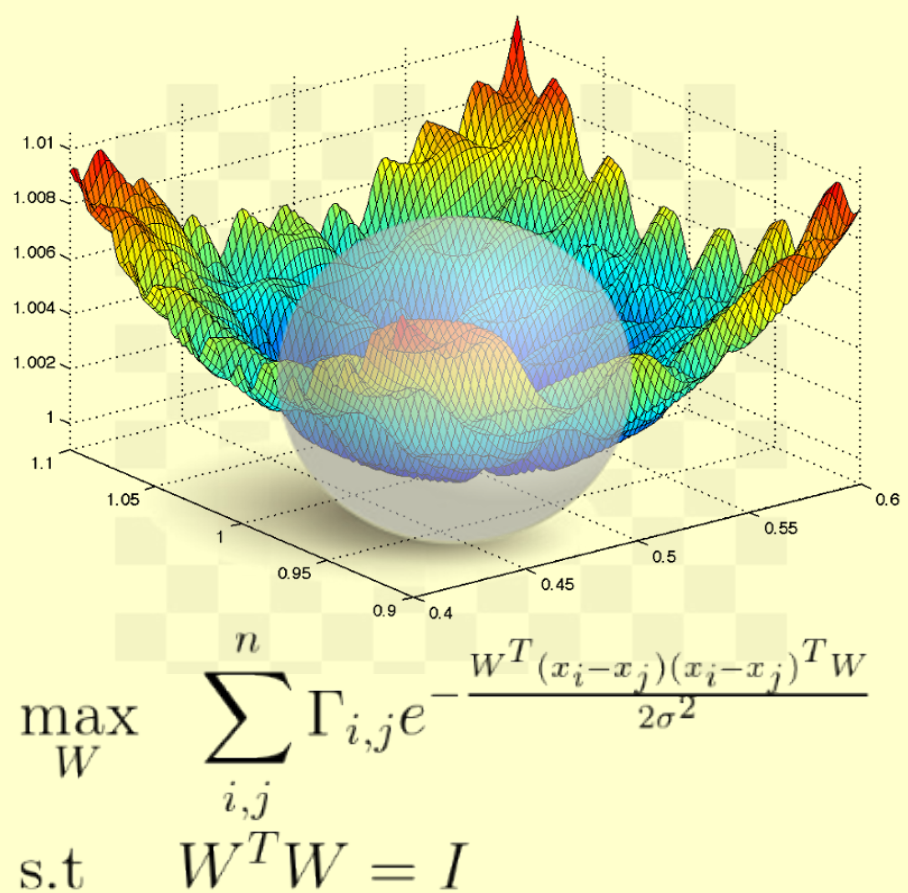
$$\text{s.t.} \quad W^T W = I, Y^T Y = I$$

**Publications that used IKDR**

- Barshan, Elnaz, et al. "Supervised principal component analysis: Visualization, classification and regression on subspaces and submanifolds." Pattern Recognition 44.7 (2011): 1357-1371.
- Masaeli, Mahdokht, Jennifer G. Dy, and Glenn M. Fung. "From transformation-based dimensionality reduction to feature selection." Proceedings of the 27th International Conference on Machine Learning (ICML-10), 2010.
- Niu, Donglin, Jennifer G. Dy, and Michael I. Jordan. "Multiple non-redundant spectral clustering views." Proceedings of the 27th international conference on machine learning (ICML-10), 2010.
- Niu, Donglin, Jennifer Dy, and Michael I. Jordan. "Dimensionality reduction for spectral clustering." Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics. 2011.
- Wu, Chieh, et al. "Iterative spectral method for alternative clustering." International Conference on Artificial Intelligence and Statistics. 2018.
- Chang, Yale, et al. "Clustering with Domain-Specific Usefulness Scores." Proceedings of the 2017 SIAM International Conference on Data Mining. Society for Industrial and Applied Mathematics, 2017.

### Why is IKDR Difficult?

**Optimizing W is highly non-convex and the solution must intersect the Stiefel Manifold**



**Existing Solutions**

Dimension Growth  
Optimization Via Stiefel Manifold.  
Optimization Via Grassmann Manifold.  
Stochastic Gradient Descent.

**Problems with Existing Solutions**

very slow  
Difficult to implement  
stuck at saddle point  
poor results

**Our Solution**

The Iterative Spectral Method ( ISM )

### Our Solution : The Iterative Spectral Method (ISM)

**We identified a special family of kernels (The ISM family) with the following properties:**

- Each kernel within the family has an associated scaled covariance matrix  $\Phi$ .
- The most dominant eigenvectors of  $\Phi$  is the solution to Eq. ( 1 ).
- The conic combination of ISM kernels is still in the ISM family.
- The conic combination of  $\Phi$ s is the associated scaled covariance matrix for the conic combination of kernels.
- If  $\Phi$  is a function of W, then  $\Phi$  can be approximated using the 2<sup>nd</sup> order Taylor series

**Formal Definition of the ISM family:**

**Definition 1.** Given  $\beta = a(x_i, x_j)^T W W^T b(x_i, x_j)$  with  $a(x_i, x_j)$  and  $b(x_i, x_j)$  as functions of  $x_i$  and  $x_j$ , any twice differentiable kernel that can be written in terms of  $f(\beta)$  while retaining its symmetric positive semi-definite property is an ISM kernel belonging to the ISM family with an associated  $\Phi$  matrix defined as

$$\Phi = \frac{1}{2} \sum_{i,j} \Gamma_{i,j} f'(\beta) A_{i,j}. \quad (6)$$

where  $A_{i,j} = b(x_i, x_j) a(x_i, x_j)^T + a(x_i, x_j) b(x_i, x_j)^T$ .

**Theorem 3.** For any kernel within the ISM family, a  $\Phi$  independent of W can be approximated with

$$\Phi \approx \text{sign}(\nabla_{\beta} f(0)) \sum_{i,j} \Gamma_{i,j} A_{i,j}. \quad (7)$$

**The ISM Algorithm:**

**Algorithm 1** ISM Algorithm

**Input :** Data X, kernel, Subspace Dimension q

**Output :** Projected subspace W

**Initialization :** Initialize  $\Phi_0$  using Table 1.

Set  $W_0$  to  $V_{\max}$  of  $\Phi_0$ .

**while**  $\|\Delta_i - \Delta_{i-1}\|_2 / \|\Delta_i\|_2 < \delta$  **do**

    Compute  $\Phi$  using Table 2

    Set  $W_k$  to  $V_{\max}$  of  $\Phi$

**end**

**Examples of Approximations of  $\Phi$ s**

Kernel	Approximation of $\Phi$ s
Linear	$\Phi_0 = X^T \Gamma X$
Squared	$\Phi_0 = X^T \mathcal{L}_{\Gamma} X$
Polynomial	$\Phi_0 = X^T \Gamma X$
Gaussian	$\Phi_0 = -X^T \mathcal{L}_{\Gamma} X$
Multiquadratic	$\Phi_0 = X^T \mathcal{L}_{\Gamma} X$

Table 1: Equations for the approximate  $\Phi$ s for the common kernels.

**Examples of of  $\Phi$ s**

Kernel	$\Phi$ Equations
Linear	$\Phi = X^T \Gamma X$
Squared	$\Phi = X^T \mathcal{L}_{\Gamma} X$
Polynomial	$\Phi = X^T \Psi X$ , $\Psi = \Gamma \odot K_{XW, p-1}$
Gaussian	$\Phi = -X^T \mathcal{L}_{\Psi} X$ , $\Psi = \Gamma \odot K_{XW}$
Multiquadratic	$\Phi = X^T \mathcal{L}_{\Psi} X$ , $\Psi = \Gamma \odot K_{XW}^{(-1)}$

Table 2: Equations for  $\Phi$ s for the common kernels.

**How K(x,x') become f(beta):**

Kernel Name	$f(\beta)$	$a(x_i, x_j)$	$b(x_i, x_j)$
Linear	$\beta$	$x_i$	$x_j$
Squared	$\beta$	$x_i - x_j$	$x_i - x_j$
Polynomial	$(\beta + c)^p$	$x_i$	$x_j$
Gaussian	$e^{-\frac{\beta}{2\sigma^2}}$	$x_i - x_j$	$x_i - x_j$
Multiquadratic	$\sqrt{\beta + c^2}$	$x_i - x_j$	$x_i - x_j$

Table 3: Converting common kernels to  $f(\beta)$ .

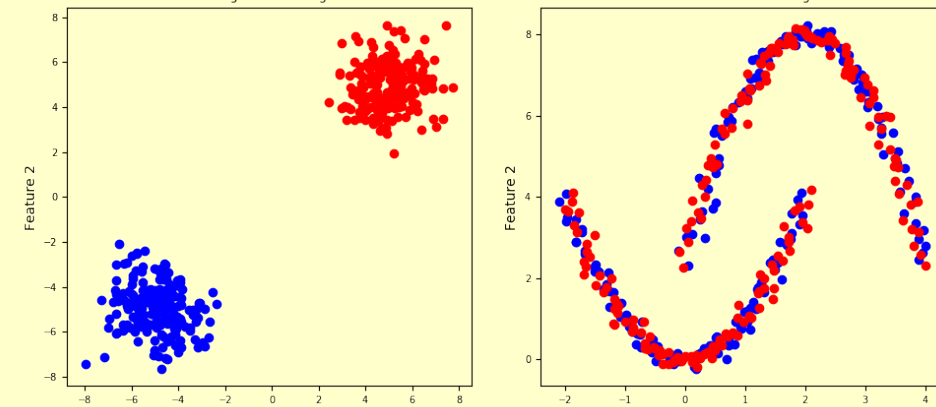
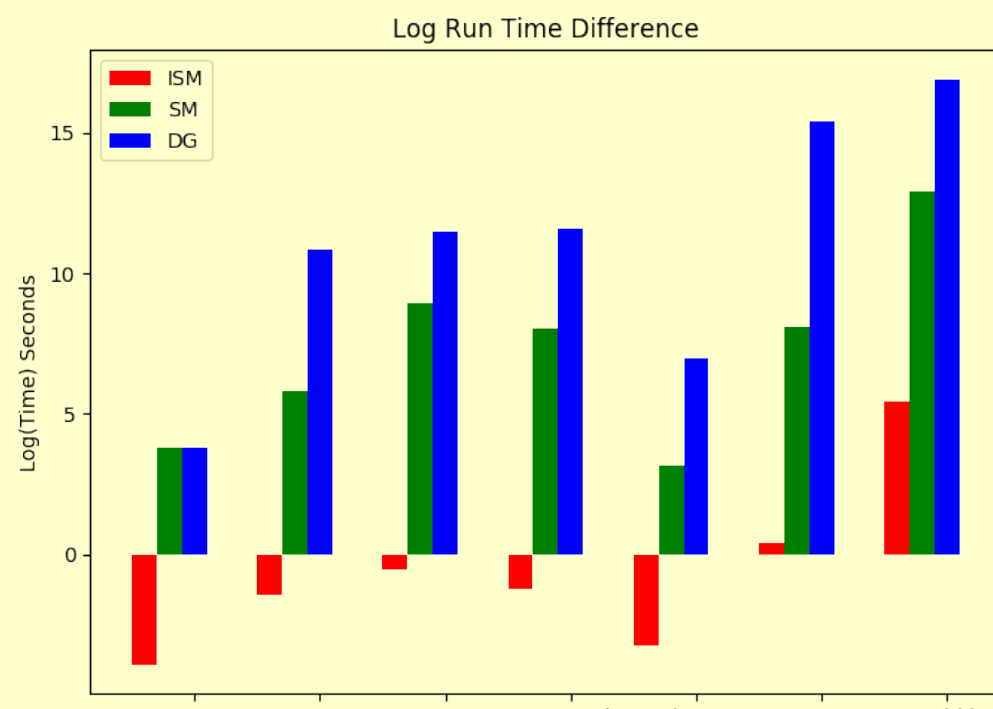
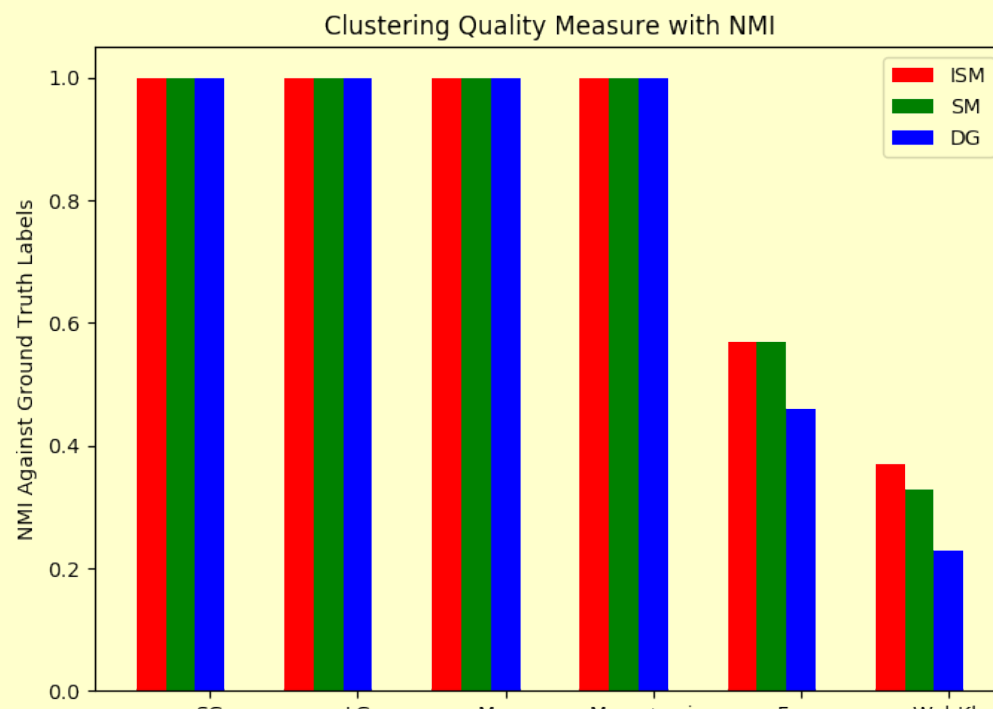
## Experimental Results

Supervised		Gaussian				polynomial			
		ISM	DG	SM	GM	ISM	DG	SM	GM
Wine	Time	<b>0.02s ± 0.01s</b>	7.9s ± 2.9s	1.7s ± 0.7s	16.8m ± 3.4s	<b>0.02s ± 0.0s</b>	13.2s ± 6.2s	14.77s ± 0.6s	16.82m ± 3.6s
	Cost	<b>-1311 ± 26</b>	-1201 ± 25	-1310 ± 26	-1307 ± 25	<b>-114608 ± 1752</b>	-112440 ± 1719	-111339 ± 1652	-108892 ± 1590
	Accuracy	<b>95.0% ± 5%</b>	93.2% ± 5.5%	<b>95% ± 4.2%</b>	<b>95% ± 6%</b>	<b>97.2% ± 3.7%</b>	93.8% ± 3.9%	96.6% ± 3.7%	96.6% ± 2.7%
Cancer	Time	<b>0.08s ± 0.0s</b>	4.5m ± 103s	17s ± 12s	17.8m ± 80s	<b>0.13s ± 0.0s</b>	4m ± 1.2m	3.3m ± 3s	17.5m ± 1.1m
	Cost	<b>-32249 ± 338</b>	-30302 ± 2297	-31996 ± 499	-30998 ± 560	<b>-1894 ± 47</b>	-1882 ± 47	-1737 ± 84	-1690 ± 108
	Accuracy	97.3% ± 0.3%	97.3% ± 0.3%	97.3% ± 0.2%	<b>97.4% ± 0.4%</b>	<b>97.4% ± 0.3%</b>	97.3% ± 0.3%	<b>97.4% ± 0.3%</b>	97.3% ± 0.3%
Face	Time	<b>0.99s ± 0.1s</b>	1.92d ± 11h	10s ± 5s	22.7m ± 18s	<b>0.7s ± 0.03s</b>	2.1d ± 13.9h	5.0m ± 5.7s	21.5m ± 9.8s
	Cost	<b>-3754 ± 31</b>	-3431 ± 32	-3749 ± 33	-771 ± 28	<b>-82407 ± 1670</b>	-78845 ± 1503	-37907 ± 15958	-3257 ± 517
	Accuracy	<b>100% ± 0%</b>	<b>100% ± 0%</b>	<b>100% ± 0%</b>	99.2% ± 0.2%	<b>100% ± 0%</b>	<b>100% ± 0%</b>	<b>100% ± 0%</b>	99.8% ± 0.2%
MNIST	Time	<b>13.8s ± 2.3s</b>	> 3d	2.5m ± 1.0s	> 3d	<b>12.1s ± 1.4s</b>	> 3d	2.1m ± 3s	> 3d
	Cost	<b>-639 ± 2.3</b>	N/A	-621 ± 5.1	N/A	<b>-639 ± 2</b>	N/A	-620 ± 5.1	N/A
	Accuracy	<b>99% ± 0%</b>	N/A	98.5% ± 0.4%	N/A	<b>99% ± 0%</b>	N/A	<b>99% ± 0%</b>	N/A
Unsupervised									
Wine	Time	<b>0.01s</b>	9.9s	0.6s	16.7m	<b>0.02s</b>	14.4s	2.9s	33.5m
	Cost	<b>-27.4</b>	-25.2	-27.3	-27.3	<b>-1600</b>	-1582	-1598	-1496
	NMI	<b>0.86</b>	<b>0.86</b>	<b>0.86</b>	<b>0.86</b>	<b>0.84</b>	<b>0.84</b>	<b>0.84</b>	0.83
Cancer	Time	<b>0.57s</b>	4.3m	3.9s	44m	<b>0.5s</b>	8.0m	8.8m	41m
	Cost	<b>-243</b>	-133	-146	-142	<b>-15804</b>	-14094	-15749	-11985
	NMI	<b>0.8</b>	0.79	<b>0.8</b>	0.79	<b>0.80</b>	<b>0.80</b>	0.79	<b>0.80</b>
Face	Time	<b>0.3s</b>	1.3d	5.3s	55.9m	<b>1.0s</b>	> 3d	22m	1.6d
	Cost	<b>-169.3</b>	-167.7	-168.9	-37	<b>-368</b>	NA	-348	-321
	NMI	0.94	<b>0.95</b>	0.93	0.89	<b>0.94</b>	N/A	0.89	0.89
MNIST	Time	<b>1.8h</b>	> 3d	1.3d	> 3d	<b>8.3m</b>	> 3d	0.9d	> 3d
	Cost	<b>-2105</b>	N/A	-2001	N/A	<b>-51358</b>	N/A	-51129	N/A
	NMI	<b>0.47</b>	N/A	0.46	N/A	<b>0.32</b>	N/A	<b>0.32</b>	N/A

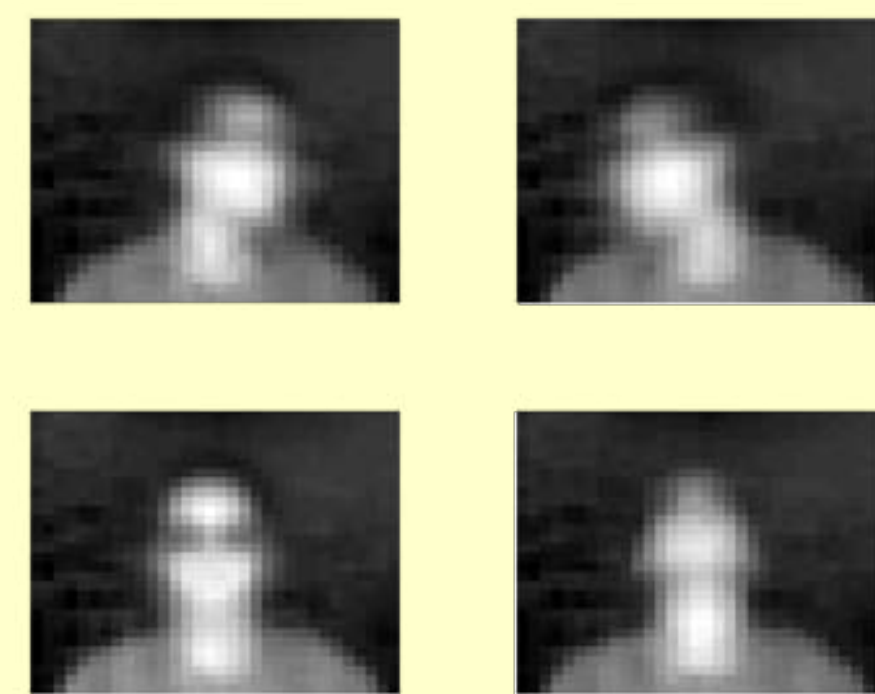
Table 4: Run-time, cost, and objective performance are recorded under supervised/unsupervised objectives. ISM is significantly faster compared to other optimization techniques while achieving lower objective cost.

		Supervised				Unsupervised		
		Linear	Squared	Multiquad	G+P	Linear	Squared	Multiquad
Wine	Time	<b>0.003s ± 0s</b>	0.01s ± 0s	0.02s ± 0.01s	0.007s ± 0s	<b>0.02s</b>	0.04s	0.06s
	Accuracy	97.2% ± 2.8%	96.6% ± 3.7%	97.2% ± 3.7%	<b>98.3% ± 2.6%</b>	0.85	0.85	<b>0.88</b>
Cancer	Time	<b>0.02s ± 0.002s</b>	0.09s ± 0.02s	0.15s ± 0.01s	0.06s ± 0.004s	<b>Time 0.23s</b>	0.5s	0.56s
	Accuracy	97.2% ± 0.3%	97.3% ± 0.04%	<b>97.4% ± 0.003%</b>	<b>97.4% ± 0.003%</b>	NMI 0.80	0.79	<b>0.84</b>
Face	Time	<b>0.2s ± 0.2s</b>	0.3s ± 0.2s	0.3s ± 0.2s	0.5s ± 0.03s	<b>Time 0.68s</b>	0.92s	3.7s
	Accuracy	97.3% ± 0.3%	97.1% ± 0.4%	97.3% ± 0.4%	<b>100% ± 0%</b>	NMI 0.93	<b>0.95</b>	0.92
MNIST	Time	<b>6.4s ± 0.4s</b>	17.4s ± 0.4s	10.6m ± 1.9m	17.6s ± 2.5s	<b>Time 3.1m</b>	4.7m	52m
	Accuracy	99.1% ± 0.1%	<b>99.3% ± 0.2%</b>	99.1% ± 0.1%	<b>99.3% ± 0.2%</b>	NMI 0.54	<b>0.54</b>	<b>0.54</b>

Table 5: Run-time and objective performance are recorded across several kernels within the ISM family. It confirms the usage of  $\Phi$  or linear combination of  $\Phi$  in place of kernels.



(a) Identity View (Mean Images)



(b) Pose View (Mean Images)

## ISM's Theoretical Foundation

**Theorem 1:**

Given a full rank  $\Phi$  with an eigengap as defined by Eq. (80), a fixed point  $W^*$  of algorithm 1 satisfies the 2<sup>nd</sup> order necessary condition using any ISM Kernel.

$$\left( \min_i \bar{\Lambda}_i - \max_j \Lambda_j \right) \geq C. \quad (80)$$

**Theorem 2:**

A sequence of subspaces generated by Algorithm 1 contains a converging subsequence.

**Theorem 3:**

For any kernel within the ISM family, a  $\Phi$  Independent of W can be approximated with

$$\Phi \approx \text{sign}(\nabla_{\beta} f(0)) \sum_{i,j} \Gamma_{i,j} A_{i,j}.$$

**Proposition 1:**

Any conic combination of ISM kernels is still an ISM kernel.

**Corollary 1:**

The  $\Phi$  matrix associated with a conic combination of kernels is the conic combination of  $\Phi$ s associated with each individual kernel.

Acknowledgments: This work was made possible by the National Science Foundation (NSF IIS-1546428). The PROTECT data is supported by the National Institute of Environmental Health Sciences (P42ES017198).